

SVD and Low-rank Approximation

Lecture 23

April 18, 2019

Singular Value Decomposition (SVD)

Let A be a $m \times n$ real-valued matrix

- a_i denotes vector corresponding to row i
- m rows. think of each row as a data point in \mathbb{R}^n
- Data applications: $m \gg n$
- Other notation: A is a $n \times d$ matrix.

Singular Value Decomposition (SVD)

Let A be a $m \times n$ real-valued matrix

- a_i denotes vector corresponding to row i
- m rows. think of each row as a data point in \mathbb{R}^n
- Data applications: $m \gg n$
- Other notation: A is a $n \times d$ matrix.

SVD theorem: A can be written as UDV^T where

- V is a $n \times n$ orthonormal matrix
- D is a $m \times n$ diagonal matrix with $\leq \min\{m, n\}$ non-zeroes called the singular values of A
- U is a $m \times m$ orthonormal matrix

SVD

Let $d = \min\{m, n\}$.

- u_1, u_2, \dots, u_m columns of U , left singular vectors of A
- v_1, v_2, \dots, v_n columns of V (rows of V^T) right singular vectors of A
- $\sigma_1 \geq \sigma_2 \geq \dots, \geq \sigma_d$ are singular values where $d = \min\{m, n\}$. And $\sigma_i = D_{i,i}$

$$A = \sum_{i=1}^d \sigma_i u_i v_i^T$$

Let $d = \min\{m, n\}$.

- u_1, u_2, \dots, u_m columns of U , left singular vectors of A
- v_1, v_2, \dots, v_n columns of V (rows of V^T) right singular vectors of A
- $\sigma_1 \geq \sigma_2 \geq \dots, \geq \sigma_d$ are singular values where $d = \min\{m, n\}$. And $\sigma_i = D_{i,i}$

$$A = \sum_{i=1}^d \sigma_i u_i v_i^T$$

We can in fact restrict attention to r the rank of A .

$$A = \sum_{i=1}^r \sigma_i u_i v_i^T$$

Interpreting \mathbf{A} as a linear operator $\mathbf{A} : \mathbb{R}^n \rightarrow \mathbb{R}^m$

- Columns of \mathbf{V} is an orthonormal basis and hence $\mathbf{V}^T \mathbf{x}$ for $\mathbf{x} \in \mathbb{R}^n$ expresses \mathbf{x} in the \mathbf{V} basis. Note that $\mathbf{V}^T \mathbf{x}$ is a rigid transformation (does not change length of \mathbf{x}).
- Let $\mathbf{y} = \mathbf{V}^T \mathbf{x}$. \mathbf{D} is a diagonal matrix which only stretches \mathbf{y} along the coordinate axes. Also adjusts dimension to go from n to m with right number of zeroes.
- Let $\mathbf{z} = \mathbf{D} \mathbf{y}$. Then $\mathbf{U} \mathbf{z}$ is a rigid transformation that expresses \mathbf{z} in the basis corresponding to rows of \mathbf{U} .

Thus any linear operator can be broken up into a sequence of three simpler/basic type of transformations

Low rank approximation property of SVD

Question: Given $\mathbf{A} \in \mathbb{R}^{m \times n}$ and integer k find a matrix \mathbf{B} of rank at most k such that $\|\mathbf{A} - \mathbf{B}\|$ is minimized

Low rank approximation property of SVD

Question: Given $A \in \mathbb{R}^{m \times n}$ and integer k find a matrix B of rank at most k such that $\|A - B\|$ is minimized

Fact: For Frobenius norm optimum for all k is captured by SVD.

That is, $A_k = \sum_{i=1}^k \sigma_i u_i v_i^T$ is the best rank k approximation to A

$$\|A - A_k\|_F = \min_{B: \text{rank}(B) \leq k} \|A - B\|_F$$

Low rank approximation property of SVD

Question: Given $\mathbf{A} \in \mathbb{R}^{m \times n}$ and integer k find a matrix \mathbf{B} of rank at most k such that $\|\mathbf{A} - \mathbf{B}\|$ is minimized

Fact: For Frobenius norm optimum for all k is captured by SVD.

That is, $\mathbf{A}_k = \sum_{i=1}^k \sigma_i \mathbf{u}_i \mathbf{v}_i^T$ is the best rank k approximation to \mathbf{A}

$$\|\mathbf{A} - \mathbf{A}_k\|_F = \min_{\mathbf{B}: \text{rank}(\mathbf{B}) \leq k} \|\mathbf{A} - \mathbf{B}\|_F$$

Why this magic? Frobenius norm and basic properties of vector projections

Geometric meaning

Consider $k = 1$. What is the best rank 1 matrix B that minimizes $\|A - B\|_F$

Since B is rank 1 , $B = uv^T$ where $v \in \mathbb{R}^n$ and $u \in \mathbb{R}^m$
Wlog v is a unit vector

Geometric meaning

Consider $k = 1$. What is the best rank **1** matrix B that minimizes $\|A - B\|_F$

Since B is rank **1**, $B = uv^T$ where $v \in \mathbb{R}^n$ and $u \in \mathbb{R}^m$
Wlog v is a unit vector

$$\|A - uv^T\|_F^2 = \sum_{i=1}^m \|a_i - u(i)v\|^2$$

Geometric meaning

Consider $k = 1$. What is the best rank **1** matrix B that minimizes $\|A - B\|_F$

Since B is rank **1**, $B = uv^T$ where $v \in \mathbb{R}^n$ and $u \in \mathbb{R}^m$
Wlog v is a unit vector

$$\|A - uv^T\|_F^2 = \sum_{i=1}^m \|a_i - u(i)v\|^2$$

If we know v then best u to minimize above is determined. Why?

Geometric meaning

Consider $k = 1$. What is the best rank **1** matrix B that minimizes $\|A - B\|_F$

Since B is rank **1**, $B = uv^T$ where $v \in \mathbb{R}^n$ and $u \in \mathbb{R}^m$
Wlog v is a unit vector

$$\|A - uv^T\|_F^2 = \sum_{i=1}^m \|a_i - u(i)v\|^2$$

If we know v then best u to minimize above is determined. Why?
For fixed v , $u(i) = \langle a_i, v \rangle$

Geometric meaning

Consider $k = 1$. What is the best rank 1 matrix B that minimizes $\|A - B\|_F$

Since B is rank 1 , $B = uv^T$ where $v \in \mathbb{R}^n$ and $u \in \mathbb{R}^m$
Wlog v is a unit vector

$$\|A - uv^T\|_F^2 = \sum_{i=1}^m \|a_i - u(i)v\|^2$$

If we know v then best u to minimize above is determined. Why?

For fixed v , $u(i) = \langle a_i, v \rangle$

$\|a_i - \langle a_i, v \rangle v\|_2$ is distance of a_i from line described by v .

Geometric meaning

What is the best rank **1** matrix B that minimizes $\|A - B\|_F$

It is to find unit vector/direction \mathbf{v} to minimize

$$\sum_{i=1}^m \|\mathbf{a}_i - \langle \mathbf{a}_i, \mathbf{v} \rangle \mathbf{v}\|^2$$

which is same as finding unit vector \mathbf{v} to maximize

$$\sum_{i=1}^m \langle \mathbf{a}_i, \mathbf{v} \rangle^2$$

Geometric meaning

What is the best rank **1** matrix B that minimizes $\|A - B\|_F$

It is to find unit vector/direction \mathbf{v} to minimize

$$\sum_{i=1}^m \|a_i - \langle a_i, \mathbf{v} \rangle \mathbf{v}\|^2$$

which is same as finding unit vector \mathbf{v} to maximize

$$\sum_{i=1}^m \langle a_i, \mathbf{v} \rangle^2$$

How to find best \mathbf{v} ? Not obvious: we will come to it a bit later

Best rank two approximation

Consider $k = 2$. What is the best rank 2 matrix B that minimizes $\|A - B\|_F$

Since B has rank 2 we can assume without loss of generality that $B = u_1 v_1^T + u_2 v_2^T$ where v_1, v_2 are orthogonal unit vectors (span a space of dimension 2)

Best rank two approximation

Consider $k = 2$. What is the best rank **2** matrix B that minimizes $\|A - B\|_F$

Since B has rank **2** we can assume without loss of generality that $B = u_1 v_1^T + u_2 v_2^T$ where v_1, v_2 are orthogonal unit vectors (span a space of dimension **2**)

Minimizing $\|A - B\|_F^2$ is same as finding orthogonal vectors v_1, v_2 to maximize

$$\sum_{i=1}^m (\langle a_i, v_1 \rangle^2 + \langle a_i, v_2 \rangle^2)$$

in other words the best fit **2**-dimensional space

Greedy algorithm

- Find \mathbf{v}_1 as the best rank **1** approximation. That is $\mathbf{v}_1 = \arg \max_{\mathbf{v}, \|\mathbf{v}\|_2=1} \sum_{i=1}^m \langle \mathbf{a}_i, \mathbf{v} \rangle^2$
- For \mathbf{v}_2 solve $\arg \max_{\mathbf{v} \perp \mathbf{v}_1, \|\mathbf{v}\|_2=1} \sum_{i=1}^m \langle \mathbf{a}_i, \mathbf{v} \rangle^2$.

Alternatively: let $\mathbf{a}'_i = \mathbf{a}_i - \langle \mathbf{a}_i, \mathbf{v}_1 \rangle \mathbf{v}_1$. Let $\mathbf{v}_2 = \arg \max_{\mathbf{v}, \|\mathbf{v}\|_2=1} \sum_{i=1}^m \langle \mathbf{a}'_i, \mathbf{v} \rangle^2$

Greedy algorithm

- Find \mathbf{v}_1 as the best rank **1** approximation. That is $\mathbf{v}_1 = \mathbf{arg\ max}_{\mathbf{v}, \|\mathbf{v}\|_2=1} \sum_{i=1}^m \langle \mathbf{a}_i, \mathbf{v} \rangle^2$
- For \mathbf{v}_2 solve $\mathbf{arg\ max}_{\mathbf{v} \perp \mathbf{v}_1, \|\mathbf{v}\|_2=1} \sum_{i=1}^m \langle \mathbf{a}_i, \mathbf{v} \rangle^2$.

Alternatively: let $\mathbf{a}'_i = \mathbf{a}_i - \langle \mathbf{a}_i, \mathbf{v}_1 \rangle \mathbf{v}_1$. Let $\mathbf{v}_2 = \mathbf{arg\ max}_{\mathbf{v}, \|\mathbf{v}\|_2=1} \sum_{i=1}^m \langle \mathbf{a}'_i, \mathbf{v} \rangle^2$

Greedy algorithm works!

Greedy algorithm correctness

Proof that Greedy works for $k = 1$.

Suppose w_1, w_2 are orthogonal unit vectors that form the best fit 2-d space. Let H be the space spanned by w_1, w_2 .

Suffices to prove that

$$\sum_{i=1}^m (\langle a_i, v_1 \rangle^2 + \langle a_i, v_2 \rangle^2) \geq \sum_{i=1}^m (\langle a_i, w_1 \rangle^2 + \langle a_i, w_2 \rangle^2)$$

Greedy algorithm correctness

Proof that Greedy works for $k = 1$.

Suppose w_1, w_2 are orthogonal unit vectors that form the best fit 2-d space. Let H be the space spanned by w_1, w_2 .

Suffices to prove that

$$\sum_{i=1}^m (\langle a_i, v_1 \rangle^2 + \langle a_i, v_2 \rangle^2) \geq \sum_{i=1}^m (\langle a_i, w_1 \rangle^2 + \langle a_i, w_2 \rangle^2)$$

If $v_1 \subset H$ then done because we can assume wlog that $w_1 = v_1$ and v_2 is at least as good as w_2 .

Greedy algorithm correctness

Suppose $v_1 \notin H$. Let v_1' be projection of v_1 onto H and $v_1'' = v_1 - v_1'$ be the component of v_1 orthogonal to H .

Greedy algorithm correctness

Suppose $\mathbf{v}_1 \notin H$. Let \mathbf{v}'_1 be projection of \mathbf{v}_1 onto H and $\mathbf{v}''_1 = \mathbf{v}_1 - \mathbf{v}'_1$ be the component of \mathbf{v}_1 orthogonal to H . Note that $\|\mathbf{v}'_1\|^2 + \|\mathbf{v}''_1\|_2^2 = \|\mathbf{v}_1\|_2^2 = \mathbf{1}$.

Wlog we can assume by rotation that $\mathbf{w}_1 = \frac{1}{\|\mathbf{v}'_1\|_2} \mathbf{v}'_1$ and \mathbf{w}_2 is orthogonal to \mathbf{v}'_1 . Hence \mathbf{w}_2 is orthogonal to \mathbf{v}_1 .

Greedy algorithm correctness

Suppose $v_1 \notin H$. Let v_1' be projection of v_1 onto H and $v_1'' = v_1 - v_1'$ be the component of v_1 orthogonal to H . Note that $\|v_1'\|^2 + \|v_1''\|_2^2 = \|v_1\|_2^2 = 1$.

Wlog we can assume by rotation that $w_1 = \frac{1}{\|v_1'\|_2} v_1'$ and w_2 is orthogonal to v_1' . Hence w_2 is orthogonal to v_1 .

Therefore v_2 is at least as good as w_2 , and v_1 is at least as good as w_1 which implies the desired claim.

Greedy algorithm for general k

- Find \mathbf{v}_1 as the best rank **1** approximation. That is $\mathbf{v}_1 = \mathbf{arg\ max}_{\mathbf{v}, \|\mathbf{v}\|_2=1} \sum_{i=1}^m \langle \mathbf{a}_i, \mathbf{v} \rangle^2$
- For \mathbf{v}_k solve $\mathbf{arg\ max}_{\mathbf{v} \perp \mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_{k-1}, \|\mathbf{v}\|_2=1} \sum_{i=1}^k \langle \mathbf{a}_i, \mathbf{v} \rangle^2$ which is same as solving $k = \mathbf{1}$ with vectors $\mathbf{a}'_1, \mathbf{a}'_2, \dots, \mathbf{a}'_m$ that are residuals. That is $\mathbf{a}'_i = \mathbf{a}_i - \sum_{j=1}^{k-1} \langle \mathbf{a}_i, \mathbf{v}_j \rangle \mathbf{v}_j$

Proof of correctness is via induction and is a straight forward generalization of the proof for $k = 2$

Summarizing

$$\sigma_j^2 = \sum_{i=1}^m \langle \mathbf{a}_i, \mathbf{v}_j \rangle^2$$

By greedy construction $\sigma_1 \geq \sigma_2 \dots$,

Let r be the (row) rank of \mathbf{A} . $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r$ span the row space of \mathbf{A} and $\sigma_j = 0$ for $j > r$

\mathbf{u}_1 determined by \mathbf{v}_1 and \mathbf{u}_2 determined by $\mathbf{v}_1, \mathbf{v}_2$ and so on. Can show that they are orthogonal.

$$\mathbf{A} = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T$$

Power method

Thus SVD relies on being able to solve $k = 1$ case

Given m vectors $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_m \in \mathbb{R}^n$ solve

$$\max_{\mathbf{v} \in \mathbb{R}^n, \|\mathbf{v}\|_2=1} \langle \mathbf{a}_i, \mathbf{v} \rangle^2$$

How do we solve the above problem?

Let $\mathbf{B} = \mathbf{A}^T \mathbf{A}$ Then

$$\begin{aligned} \mathbf{B} &= \left(\sum_{i=1}^m \sigma_i \mathbf{v}_i \mathbf{u}_i^T \right) \left(\sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T \right) \\ &= \sum_{i=1}^r \sigma_i^2 \mathbf{v}_i \mathbf{v}_i^T \end{aligned}$$

Power method continued

Let $B = A^T A$ Then

$$\begin{aligned} B^2 &= \left(\sum_{i=1}^r \sigma_i^2 v_i v_i^T \right) \left(\sum_{i=1}^r \sigma_i^2 v_i v_i^T \right) \\ &= \sum_{i=1}^r \sigma_i^4 v_i v_i^T. \end{aligned}$$

More generally

$$B^k = \sum_{i=1}^r \sigma_i^k v_i v_i^T$$

Power method continued

Let $B = A^T A$ Then

$$\begin{aligned} B^2 &= \left(\sum_{i=1}^r \sigma_i^2 v_i v_i^T \right) \left(\sum_{i=1}^r \sigma_i^2 v_i v_i^T \right) \\ &= \sum_{i=1}^r \sigma_i^4 v_i v_i^T. \end{aligned}$$

More generally

$$B^k = \sum_{i=1}^r \sigma_i^k v_i v_i^T$$

If $\sigma_1 > \sigma_2$ then B^k converges to $\sigma_1^k v_1 v_1^T$ and we can identify v_1 from B^k . But expensive to compute B^k

Power method continued

Pick a random (unit) vector $\mathbf{x} \in \mathbb{R}^n$. Then $\mathbf{x} = \sum_{i=1}^n \lambda_i \mathbf{v}_i$ since $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$ is a basis for \mathbb{R}^n .

$$\mathbf{B}^k \mathbf{x} = \left(\sum_{i=1}^r \sigma_i^k \mathbf{v}_i \mathbf{v}_i^T \right) \left(\sum_{i=1}^d \lambda_i \mathbf{v}_i \right) \rightarrow \sigma_1^{2k} \lambda_1 \mathbf{v}_1$$

Can obtain \mathbf{v}_1 by normalizing $\mathbf{B}^k \mathbf{x}$ to a unit vector.

Computing $\mathbf{B}^k \mathbf{x}$ is easier via a series of matrix vector multiplications

Power method continued

Pick a random (unit) vector $x \in \mathbb{R}^n$. Then $x = \sum_{i=1}^n \lambda_i v_i$ since v_1, v_2, \dots, v_n is a basis for \mathbb{R}^n .

$$B^k x = \left(\sum_{i=1}^r \sigma_i^k v_i v_i^T \right) \left(\sum_{i=1}^d \lambda_i v_i \right) \rightarrow \sigma_1^{2k} \lambda_1 v_1$$

Can obtain v_1 by normalizing $B^k x$ to a unit vector.

Computing $B^k x$ is easier via a series of matrix vector multiplications

Why random x ?

What if $\sigma_1 \simeq \sigma_2$? Power method still works. See references.

Linear least square/Regression and SVD

Linear least squares: Given $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$ find \mathbf{x} to minimize $\|\mathbf{Ax} - \mathbf{b}\|_2$.

Interesting when $m > n$ the over constrained case when there is no solution to $\mathbf{Ax} = \mathbf{b}$ and want to find best fit.

Linear least square/Regression and SVD

Linear least squares: Given $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$ find \mathbf{x} to minimize $\|\mathbf{Ax} - \mathbf{b}\|_2$.

Interesting when $m > n$ the over constrained case when there is no solution to $\mathbf{Ax} = \mathbf{b}$ and want to find best fit.

Geometrically \mathbf{Ax} is a linear combination of columns of \mathbf{A} . Hence we are asking what is the vector \mathbf{z} in the column space of \mathbf{A} that is closest to vector \mathbf{b} in ℓ_2 norm.

Linear least square/Regression and SVD

Linear least squares: Given $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$ find \mathbf{x} to minimize $\|\mathbf{Ax} - \mathbf{b}\|_2$.

Interesting when $m > n$ the over constrained case when there is no solution to $\mathbf{Ax} = \mathbf{b}$ and want to find best fit.

Geometrically \mathbf{Ax} is a linear combination of columns of \mathbf{A} . Hence we are asking what is the vector \mathbf{z} in the column space of \mathbf{A} that is closest to vector \mathbf{b} in ℓ_2 norm.

Closest vector to \mathbf{b} is the projection of \mathbf{b} into the column space of \mathbf{A} so it is “obvious” geometrically. How do we find it?

Linear least square/Regression and SVD

Linear least squares: Given $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$ find \mathbf{x} to minimize $\|\mathbf{Ax} - \mathbf{b}\|_2$.

Interesting when $m > n$ the over constrained case when there is no solution to $\mathbf{Ax} = \mathbf{b}$ and want to find best fit.

Geometrically \mathbf{Ax} is a linear combination of columns of \mathbf{A} . Hence we are asking what is the vector \mathbf{z} in the column space of \mathbf{A} that is closest to vector \mathbf{b} in ℓ_2 norm.

Closest vector to \mathbf{b} is the projection of \mathbf{b} into the column space of \mathbf{A} so it is “obvious” geometrically. How do we find it? Find an orthonormal basis $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_r$ for the columns of \mathbf{A} . Compute projection \mathbf{b}' as $\mathbf{b}' = \sum_{j=1}^r \langle \mathbf{b}, \mathbf{z}_j \rangle \mathbf{z}_j$ and output answer as $\|\mathbf{b} - \mathbf{b}'\|_2$.

Linear least square/Regression and SVD

Linear least squares: Given $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$ find \mathbf{x} to minimize $\|\mathbf{Ax} - \mathbf{b}\|_2$.

Closest vector to \mathbf{b} is the projection of \mathbf{b} into the column space of \mathbf{A} so it is “obvious” geometrically. Find an orthonormal basis $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_r$ for the columns of \mathbf{A} . Compute projection \mathbf{b}' as $\mathbf{b}' = \sum_{j=1}^r \langle \mathbf{b}, \mathbf{z}_j \rangle \mathbf{z}_j$ and output answer as $\|\mathbf{b} - \mathbf{b}'\|_2$.

Finding the basis is the expensive part. Recall SVD gives $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r$ which form a basis for the *row* space of \mathbf{A} but then $\mathbf{u}_1^T, \mathbf{u}_2^T, \dots, \mathbf{u}_m^T$ form a basis for the *column* space of \mathbf{A} . Hence SVD gives us all the information to find \mathbf{b}' . In fact we have

$$\min_{\mathbf{x}} \|\mathbf{Ax} - \mathbf{b}\|_2^2 = \sum_{i=r+1}^m \langle \mathbf{u}_i^T, \mathbf{b} \rangle^2$$